



# Certifying and reasoning on cost annotations in C programs

Nicolas Ayache, Roberto Amadio, Yann Régis-Gianas

## ► To cite this version:

Nicolas Ayache, Roberto Amadio, Yann Régis-Gianas. Certifying and reasoning on cost annotations in C programs. FMICS 2012 - 17th International Workshop on Formal Methods for Industrial Critical Systems, Aug 2012, Paris, France. hal-00702665v2

**HAL Id: hal-00702665**

**<https://inria.hal.science/hal-00702665v2>**

Submitted on 31 May 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Certifying and reasoning on cost annotations in C programs

Nicolas Ayache<sup>1,2</sup> Roberto M. Amadio<sup>1</sup> Yann Régis-Gianas<sup>1,2</sup>

<sup>1</sup> Université Paris Diderot (UMR-CNRS 7126)

<sup>2</sup> INRIA (Team  $\pi r^2$ )

**Abstract** We present a so-called labelling method to enrich a compiler in order to turn it into a “cost annotating compiler”, that is, a compiler which can *lift* pieces of information on the execution cost of the object code as cost annotations on the source code. These cost annotations characterize the execution costs of code fragments of constant complexity. The first contribution of this paper is a proof methodology that extends standard simulation proofs of compiler correctness to ensure that the cost annotations on the source code are sound and precise with respect to an execution cost model of the object code.

As a second contribution, we demonstrate that our label-based instrumentation is scalable because it consists in a modular extension of the compilation chain. To that end, we report our successful experience in implementing and testing the labelling approach on top of a prototype compiler written in `ocaml` for (a large fragment of) the C language.

As a third and last contribution, we provide evidence for the usability of the generated cost annotations as a mean to reason on the concrete complexity of programs written in C. For this purpose, we present a FRAMA-C plugin that uses our cost annotating compiler to automatically infer trustworthy logic assertions about the concrete worst case execution cost of programs written in a fragment of the C language. These logic assertions are synthetic in the sense that they characterize the cost of executing the entire program, not only constant-time fragments. (These bounds may depend on the size of the input data.) We report our experimentations on some C programs, especially programs generated by a compiler for the synchronous programming language LUSTRE used in critical embedded software.

## 1 Introduction

The formal description and certification of software components is reaching a certain level of maturity with impressing case studies ranging from compilers to kernels of operating systems. A well-documented example is the proof of functional correctness of a moderately optimizing compiler from a large subset of the C language to a typical assembly language of the kind used in embedded systems [10].

In the framework of the *Certified Complexity* (CerCo) project<sup>1</sup> [3], we aim to refine this line of work by focusing on the issue of the *execution cost* of

---

<sup>1</sup> CerCo project <http://cerco.cs.unibo.it>

the compiled code. Specifically, we aim to build a formally verified C compiler that given a source program produces automatically a functionally equivalent object code plus an annotation of the source code which is a sound and precise description of the execution cost of the object code.

We target in particular the kind of C programs produced for embedded applications; these programs are eventually compiled to binaries executable on specific processors. The current state of the art in commercial products such as Scade<sup>2</sup> [7] is that the *reaction time* of the program is estimated by means of abstract interpretation methods (such as those developed by AbsInt<sup>3</sup> [6]) that operate on the binaries. These methods rely on a specific knowledge of the architecture of the processor and may require explicit (and uncertified) annotations of the binaries to determine the number of times a loop is iterated (see, *e.g.*, [13] for a survey of the state of the art).

In this context, our aim is to produce a mechanically verified compiler which can *lift* in a provably correct way the pieces of information on the execution cost of the binary code to cost annotations on the source C code. Then the produced cost annotations are manipulated with the Frama – C<sup>4</sup> [4] automatic tool to infer synthetic cost annotations. We stress that the practical relevance of the proposed approach depends on the possibility of obtaining accurate information on the execution cost of relatively short sequences of binary instructions. This seems beyond the scope of current Worst-Case Execution Time (WCET) tools such as AbsInt or Chronos<sup>5</sup> which do not support a *compositional* analysis of WCET. For this reason, we focus on processors with a simple architecture for which manufacturers can provide accurate information on the execution cost of the binary instructions. In particular, our experiments are based on the 8051 [9]<sup>6</sup>. This is a widely popular 8-bits processor developed by Intel for use in embedded systems with no cache and no pipeline. An important characteristic of the processor is that its cost model is ‘additive’: the cost of a sequence of instructions is exactly the sum of the costs of each instruction.

The rest of the paper is organized as follows. Section 2 describes the labelling approach and its formal application to a toy compiler. Appendix A gives standard definitions for the toy compiler and sketches the proofs. A formal and browsable Coq development composed of 1 *Kloc* of specifications and 3.5 *Kloc* of proofs is available at <http://www.pps.univ-paris-diderot.fr/cerco>. Section 3 reports our experience in implementing and testing the labelling approach for a compiler from C to 8051 binaries. The CerCo compiler is composed of 30 *Kloc* of ocaml code; it can be both downloaded and tested as a web application at the URL above. More details are available in Appendix B. Section 4 introduces the automatic Cost tool that starting from the cost annotations produces

<sup>2</sup> Esterel Technologies. <http://www.esterel-technologies.com>.

<sup>3</sup> AbsInt Angewandte Informatik. <http://www.absint.com/>.

<sup>4</sup> Frama – C software analyzers. <http://frama-c.com/>.

<sup>5</sup> Chronos tool. [www.comp.nus.edu.sg/~rpembed/chronos](http://www.comp.nus.edu.sg/~rpembed/chronos).

<sup>6</sup> The recently proposed ARM Cortex M series would be another obvious candidate.

certified synthetic cost bounds. This is a `Frama - C` plug-in composed of 5 *Kloc* of `ocaml` code also available at the URL above.

## 2 A “labelling” method for cost annotating compilation

In this section, we explain in general terms the so-called “labelling” method to turn a compiler into a cost annotating compiler while minimizing the impact of this extension on the proof of the semantic preservation. Then to make our purpose technically precise, we apply the method to a toy compiler.

### 2.1 Overview

As a first step, we need a clear and flexible picture of: (i) the meaning of cost annotations, (ii) the method to provide them being sound and precise, and (iii) the way such proofs can be composed. The execution cost of the source programs we are interested in depends on their control structure. Typically, the source programs are composed of mutually recursive procedures and loops and their execution cost depends, up to some multiplicative constant, on the number of times procedure calls and loop iterations are performed. Producing a *cost annotation* of a source program amounts to:

- enrich the program with a collection of *global cost variables* to measure resource consumption (time, stack size, heap size,...)
- inject suitable code at some critical points (procedures, loops,...) to keep track of the execution cost.

Thus, producing a cost-annotation of a source program  $P$  amounts to build an *annotated program*  $An(P)$  which behaves as  $P$  while self-monitoring its execution cost. In particular, if we do *not* observe the cost variables then we expect the annotated program  $An(P)$  to be functionally equivalent to  $P$ . Notice that in the proposed approach an annotated program is a program in the source language. Therefore, the meaning of the cost annotations is automatically defined by the semantics of the source language and tools developed to reason on the source programs can be directly applied to the annotated programs too. Finally, notice that the annotated program  $An(P)$  is *only* meant to *reason* on the execution cost of the unannotated program  $P$ . Therefore,  $An(P)$  will never be compiled, nor executed. Hence, the influence of cost annotations on the WCET is off-topic.

*Soundness and precision of cost annotations* Suppose we have a functionally correct compiler  $\mathcal{C}$  that associates with a program  $P$  in the source language a program  $\mathcal{C}(P)$  in the object language. Further suppose we have some obvious way of defining the execution cost of an object code. For instance, we have a good estimate of the number of cycles needed for the execution of each instruction of the object code. Now, the annotation of the source program  $An(P)$  is *sound* if its prediction of the execution cost is an upper bound for the ‘real’ execution cost. Moreover, we say that the annotation is *precise* with respect to the cost model if the *difference* between the predicted and real execution costs is bounded by a constant which only depends on the program.

## 2.2 The labelling approach, formally

$$\begin{array}{ccccc}
L_1 & \xleftarrow{\mathcal{I}} & L_{1,\ell} & \xrightarrow{c_1} & L_{2,\ell} & \cdots & \xrightarrow{c_k} & L_{k+1,\ell} \\
& & \uparrow \mathcal{L} & & \downarrow er_2 & & & \downarrow er_{k+1} \\
& & L_1 & \xrightarrow{c_1} & L_2 & \cdots & \xrightarrow{c_k} & L_{k+1}
\end{array}
\quad \left| \begin{array}{l}
er_{i+1} \circ \mathcal{C}_i = \mathcal{C}_i \circ er_i \\
er_1 \circ \mathcal{L} = id_{L_1} \\
An = \mathcal{I} \circ \mathcal{L}
\end{array} \right.$$

For each labelled language there is an obvious function  $er_i$  erasing all labels and producing a program in the corresponding unlabelled language. The compilation functions  $C_i$  are extended from the unlabelled to the labelled language so that they enjoy commutation with the erasure functions. Moreover, we lift

the soundness properties of the compilation functions from the unlabelled to the labelled languages and transition systems.

A *labelling*  $\mathcal{L}$  of the source language  $L_1$  is a function such that  $er_{L_1} \circ \mathcal{L}$  is the identity function. An *instrumentation*  $\mathcal{I}$  of the source labelled language  $L_{1,\ell}$  is a function replacing the labels with suitable increments of, say, a fresh global *cost* variable. Then, an *annotation*  $An$  of the source program can be derived simply as the composition of the labelling and the instrumentation functions:  $An = \mathcal{I} \circ \mathcal{L}$ .

Suppose  $s$  is some adequate representation of the state of a program. Let  $P$  be a source program. The judgement  $(P, s) \Downarrow s'$  is the big-step evaluation of  $P$  transforming state  $s$  into a state  $s'$ . Let us write  $s[v/x]$  to denote a state  $s$  in which the variable  $x$  is assigned a value  $v$ . Suppose now that its annotation satisfies the following property:

$$(An(P), s[c/cost]) \Downarrow s'[c + \delta/cost] \quad (1)$$

where  $c$  and  $\delta$  are some non-negative numbers. Then, the definition of the instrumentation and the fact that the soundness proofs of the compilation functions have been lifted to the labelled languages allows to conclude that

$$(\mathcal{C}(\mathcal{L}(P)), s[c/cost]) \Downarrow (s'[c/cost], \lambda) \quad (2)$$

where  $\mathcal{C} = \mathcal{C}_k \circ \dots \circ \mathcal{C}_1$  and  $\lambda$  is a sequence (or a multi-set) of labels whose ‘cost’ corresponds to the number  $\delta$  produced by the annotated program. Then, the commutation properties of erasure and compilation functions allows to conclude that the *erasure* of the compiled labelled code  $er_{k+1}(\mathcal{C}(\mathcal{L}(P)))$  is actually equal to the compiled code  $\mathcal{C}(P)$  we are interested in. Given this, the following question arises: under which conditions the sequence  $\lambda$ , *i.e.*, the increment  $\delta$ , is a sound and possibly precise description of the execution cost of the object code?

To answer this question, we observe that the object code we are interested in is some kind of assembly code and its control flow can be easily represented as a control flow graph. The idea is then to perform two simple checks on the control flow graph. The first check is to verify that all loops go through a labelled node. If this is the case then we can associate a finite cost with every label and prove that the cost annotations are sound. The second check amounts to verify that all paths starting from a label have the same cost. If this check is successful then we can conclude that the cost annotations are precise.

### 2.3 A toy compiler

As a first case study, we apply the labelling approach to a *toy compiler*.

The syntax of the source, intermediate and target languages is given in Figure 1. The three languages considered can be shortly described as follows: **Imp** is a very simple imperative language with pure expressions, branching and looping commands, **Vm** is an assembly-like language enriched with a stack, and **Mips** is a Mips-like assembly language [8] with registers and main memory.

The semantics of **Imp** is defined over configurations  $(S, K, s)$  where  $S$  is a statement,  $K$  is a continuation and  $s$  is a state. A *continuation*  $K$  is a list of

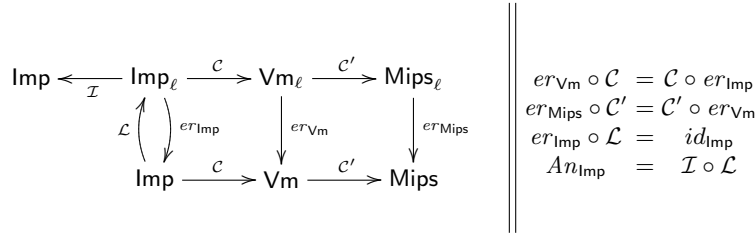
<u>Syntax for Imp</u>	<u>Syntax for Vm</u>
$id ::= x \mid y \mid \dots$	$instr_{Vm} ::= \text{cnst}(n) \mid \text{var}(n)$
$n ::= 0 \mid -1 \mid +1 \mid \dots$	$\mid \text{setvar}(n) \mid \text{add}$
$v ::= n \mid \text{true} \mid \text{false}$	$\mid \text{branch}(k) \mid \text{bge}(k) \mid \text{halt}$
$e ::= id \mid n \mid e + e$	
$b ::= e < e$	
$S ::= \text{skip} \mid id := e \mid S; S$	<u>Syntax for Mips</u>
$\mid \text{if } b \text{ then } S \text{ else } S$	$instr_{Mips} ::= \text{loadi } R, n \mid \text{load } R, l$
$\mid \text{while } b \text{ do } S$	$\mid \text{store } R, l \mid \text{add } R, R, R$
$P ::= \text{prog } S$	$\mid \text{branch } k \mid \text{bge } R, R, k \mid \text{halt}$

Figure1: Syntax definitions.

commands which terminates with a special symbol **halt**. The semantics of **Vm** is defined over stack-based machine configurations  $C \vdash (i, \sigma, s)$  where  $C$  is a program,  $i$  is a program counter,  $\sigma$  is a stack and  $s$  is a state. The semantics of **Mips** is defined over register-based machine configurations  $C \vdash (i, m)$  where  $C$  is a program,  $i$  is a program counter and  $m$  is a machine memory (with registers and main memory).

The first compilation function  $\mathcal{C}$  relies on the stack of the **Vm** language to implement expression evaluation while the second compilation function  $\mathcal{C}'$  allocates (statically) the base of the stack in the registers and the rest in main memory. This is of course a naive strategy but it suffices to expose some of the problems that arise in defining a compositional approach. The formal definitions of these compilation functions  $\mathcal{C}$  from **Imp** to **Vm** and  $\mathcal{C}'$  from **Vm** to **Mips** are standard and thus eluded. (See Appendix A for formal details about semantics and the compilation chain.)

Applying the labelling approach to this toy compiler results in the following diagram. The next sections aim at describing this diagram in details.



## 2.4 Labelled languages: Syntax and Semantics

*Syntax* The syntax of **Imp** is extended so that statements can be labelled:  $S ::= \dots \mid \ell : S$ . A new instruction **emit**( $\ell$ ) (resp. (**emit**  $\ell$ )) is introduced in the syntax of **Vm** (resp. **Mips**).

*Semantics* The small step semantics of **Imp** statements is extended as described by the following rule.

$$(\ell : S, K, s) \xrightarrow{\ell} (S, K, s)$$

We denote with  $\lambda, \lambda', \dots$  finite sequences of labels. In particular, the empty sequence is written  $\epsilon$ . We also identify an unlabelled transition with a transition labelled with  $\epsilon$ . Then, the small step reduction relation we have defined on statements becomes a *labelled transition system*. We derive a *labelled* big-step semantics as follows:  $(S, s) \Downarrow (s', \lambda)$  if  $(S, \text{halt}, s) \xrightarrow{\lambda_1} \dots \xrightarrow{\lambda_n} (\text{skip}, \text{halt}, s')$  and  $\lambda = \lambda_1 \dots \lambda_n$ .

Following the same pattern, the small step semantics of **Vm** and **Mips** are turned into a labelled transition system as follows:

$$\begin{array}{ll} C \vdash (i, \sigma, s) \xrightarrow{\ell} (i+1, \sigma, s) & \text{if } C[i] = \text{emit}(\ell) . \\ M \vdash (i, m) \xrightarrow{\ell} (i+1, m) & \text{if } M[i] = (\text{emit } \ell) . \end{array}$$

The evaluation predicate for labelled **Vm** is defined as  $(C, s) \Downarrow (s', \lambda)$  if  $C \vdash (0, \epsilon, s) \xrightarrow{\lambda_1} \dots \xrightarrow{\lambda_n} (i, \epsilon, s')$ ,  $\lambda = \lambda_1 \dots \lambda_n$  and  $C[i] = \text{halt}$ . The evaluation predicate for labelled **Mips** is defined as  $(M, m) \Downarrow (m', \lambda)$  if  $M \vdash (0, m) \xrightarrow{\lambda_1} \dots \xrightarrow{\lambda_n} (j, m')$ ,  $\lambda = \lambda_1 \dots \lambda_n$  and  $M[j] = \text{halt}$ .

## 2.5 Erasure functions

There is an obvious *erasure* function  $er_{\text{Imp}}$  from the labelled language to the unlabelled one which is the identity on expressions and boolean conditions, and traverses commands removing all labels.

The erasure function  $er_{\text{Vm}}$  amounts to remove from a **Vm** code  $C$  all the  $\text{emit}(\ell)$  instructions and recompute jumps accordingly. Specifically, let  $n(C, i, j)$  be the number of  $\text{emit}$  instructions in the interval  $[i, j]$ . Then, assuming  $C[i] = \text{branch}(k)$  we replace the offset  $k$  with an offset  $k'$  determined as follows:

$$k' = \begin{cases} k - n(C, i, i+k) & \text{if } k \geq 0 \\ k + n(C, i+1+k, i) & \text{if } k < 0 \end{cases}$$

The *erasure function*  $er_{\text{Mips}}$  is also similar to the one of **Vm** as it amounts to remove from a **Mips** code all the  $(\text{emit } \ell)$  instructions and recompute jumps accordingly. The compilation function  $C'$  is extended to  $\text{Vm}_\ell$  by simply translating  $\text{emit}(\ell)$  as  $(\text{emit } \ell)$ :

$$C'(i, C) = (\text{emit } \ell) \text{ if } C[i] = \text{emit}(\ell)$$



## 2.6 Compilation of labelled languages

The compilation function  $\mathcal{C}$  is extended to  $\text{Imp}_\ell$  by defining:

$$\mathcal{C}(\ell : b, k) = (\text{emit}(\ell)) \cdot \mathcal{C}(b, k) \quad \mathcal{C}(\ell : S) = (\text{emit}(\ell)) \cdot \mathcal{C}(S) .$$

**Proposition 1.** *For all commands  $S$  in  $\text{Imp}_\ell$ , we have that:*

- (1)  $er_{\text{Vm}}(\mathcal{C}(S)) = \mathcal{C}(er_{\text{Imp}}(S))$ .
- (2) *If  $(S, s) \Downarrow (s', \lambda)$  then  $(\mathcal{C}(S), s) \Downarrow (s', \lambda)$ .*

The following proposition relates  $\text{Vm}_\ell$  code and its compilation and it is similar to proposition 1. Here  $m \Vdash \sigma, s$  means “the low-level Mips memory  $m$  realizes the Vm stack  $\sigma$  and state  $s$ ”.

**Proposition 2.** *Let  $C$  be a  $\text{Vm}_\ell$  code. Then:*

- (1)  $er_{\text{Mips}}(\mathcal{C}'(C)) = \mathcal{C}'(er_{\text{Vm}}(C))$ .
- (2) *If  $(C, s) \Downarrow (s', \lambda)$  and  $m \Vdash \epsilon, s$  then  $(\mathcal{C}'(C), m) \Downarrow (m', \lambda)$  and  $m' \Vdash \epsilon, s'$ .*

## 2.7 Labellings and instrumentations

Assuming a function  $\kappa$  which associates an integer number with labels and a distinct variable  $cost$  which does not occur in the program  $P$  under consideration, we abbreviate with  $inc(\ell)$  the assignment  $cost := cost + \kappa(\ell)$ . Then we define the instrumentation  $\mathcal{I}$  (relative to  $\kappa$  and  $cost$ ) as follows:

$$\mathcal{I}(\ell : S) = inc(\ell); \mathcal{I}(S) .$$

The function  $\mathcal{I}$  just distributes over the other operators of the language. We extend the function  $\kappa$  on labels to sequences of labels by defining  $\kappa(\ell_1, \dots, \ell_n) = \kappa(\ell_1) + \dots + \kappa(\ell_n)$ . The instrumented  $\text{Imp}$  program relates to the labelled one as follows.

**Proposition 3.** *Let  $S$  be an  $\text{Imp}_\ell$  command. If  $(\mathcal{I}(S), s[c/cost]) \Downarrow s'[c + \delta/cost]$  then  $\exists \lambda \ \kappa(\lambda) = \delta$  and  $(S, s[c/cost]) \Downarrow (s'[c/cost], \lambda)$ .*

**Definition 1.** *A labelling is a function  $\mathcal{L}$  from an unlabelled language to the corresponding labelled one such that  $er_{\text{Imp}} \circ \mathcal{L}$  is the identity function on the  $\text{Imp}$  language.*

**Proposition 4.** *For any labelling function  $\mathcal{L}$ , and  $\text{Imp}$  program  $P$ , the following holds:*

$$er_{\text{Mips}}(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P)))) = \mathcal{C}'(\mathcal{C}(P)) . \quad (3)$$

**Proposition 5.** *Given a function  $\kappa$  for the labels and a labelling function  $\mathcal{L}$ , for all programs  $P$  of the source language if  $(\mathcal{I}(\mathcal{L}(P)), s[c/cost]) \Downarrow s'[c + \delta/cost]$  and  $m \Vdash \epsilon, s[c/cost]$  then  $(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P))), m) \Downarrow (m', \lambda)$ ,  $m' \Vdash \epsilon, s'[c/cost]$  and  $\kappa(\lambda) = \delta$ .*

## 2.8 Sound and precise labellings

With any  $\text{Mips}_\ell$  code  $M$ , we can associate a directed and rooted (control flow) graph whose nodes are the instruction positions  $\{0, \dots, |M| - 1\}$ , whose root is the node 0, and whose directed edges correspond to the possible transitions between instructions. We say that a node is labelled if it corresponds to an instruction emit  $\ell$ .

**Definition 2.** A simple path in a  $\text{Mips}_\ell$  code  $M$  is a directed finite path in the graph associated with  $M$  where the first node is labelled, the last node is the predecessor of either a labelled node or a leaf, and all the other nodes are unlabelled.

**Definition 3.** A  $\text{Mips}_\ell$  code  $M$  is soundly labelled if in the associated graph the root node 0 is labelled and there are no loops that do not go through a labelled node. Besides, we say that a soundly labelled code is precise if for every label  $\ell$  in the code, the simple paths starting from a node labelled with  $\ell$  have the same cost.

In a soundly labelled graph there are finitely many simple paths. Thus, given a soundly labelled  $\text{Mips}$  code  $M$ , we can associate with every label  $\ell$  a number  $\kappa(\ell)$  which is the maximum (estimated) cost of executing a simple path whose first node is labelled with  $\ell$ . Thus for a soundly labelled  $\text{Mips}$  code the sequence of labels associated with a computation is a significant information on the execution cost.

For an example of command which is not soundly labelled, consider  $\ell : \text{while } 0 < x \text{ do } x := x + 1$ , which when compiled, produces a loop that does not go through any label. On the other hand, for an example of a program which is not precisely labelled consider  $\ell : (\text{if } 0 < x \text{ then } x := x + 1 \text{ else skip})$ . In the compiled code, we find two simple paths associated with the label  $\ell$  whose cost will be quite different in general.

**Proposition 6.** If  $M$  is soundly (resp. precisely) labelled and  $(M, m) \Downarrow (m', \lambda)$  then the cost of the computation is bounded by  $\kappa(\lambda)$  (resp. is exactly  $\kappa(\lambda)$ ).

The next point we have to check is that there are labelling functions (of the source code) such that the compilation function does produce sound and possibly precise labelled  $\text{Mips}$  code. To discuss this point, we introduce in table 1 a labelling function  $\mathcal{L}_p$  for the  $\text{Imp}$  language. This function relies on a function “new” which is meant to return fresh labels and on an auxiliary function  $\mathcal{L}'_p$  which returns a labelled command and a binary directive  $d \in \{0, 1\}$ . If  $d = 1$  then the command that follows (if any) must be labelled.

**Proposition 7.** For all  $\text{Imp}$  programs  $P$ ,  $\mathcal{C}'(\mathcal{C}(\mathcal{L}_p(P)))$  is a soundly and precisely labelled  $\text{Mips}$  code.

Once a sound and possibly precise labelling  $\mathcal{L}$  has been designed, we can determine the cost of each label and define an instrumentation  $\mathcal{I}$  whose composition with  $\mathcal{L}$  will produce the desired cost annotation.

$\mathcal{L}_p(\text{prog } S)$	$= \text{prog } \mathcal{L}_p(S)$
$\mathcal{L}_p(S)$	$= \text{let } \ell = \text{new}, (S', d) = \mathcal{L}'_p(S) \text{ in } \ell : S'$
$\mathcal{L}'_p(S)$	$= (S, 0) \quad \text{if } S = \text{skip} \text{ or } S = (x := e)$
$\mathcal{L}'_p(\text{if } b \text{ then } S_1 \text{ else } S_2)$	$= (\text{if } b \text{ then } \mathcal{L}_p(S_1) \text{ else } \mathcal{L}_p(S_2), 1)$
$\mathcal{L}'_p(\text{while } b \text{ do } S)$	$= (\text{while } b \text{ do } \mathcal{L}_p(S), 1)$
$\mathcal{L}'_p(S_1; S_2)$	$= \text{let } (S'_1, d_1) = \mathcal{L}'_p(S_1), (S'_2, d_2) = \mathcal{L}'_p(S_2) \text{ in}$ $\quad \text{case } d_1$ $\quad 0 : (S'_1; S'_2, d_2)$ $\quad 1 : \text{let } \ell = \text{new in } (S'_1; \ell : S'_2, d_2)$

Table1: A labelling for the `Imp` language

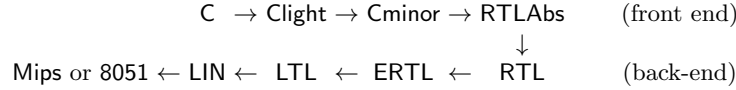
**Definition 4.** *Given a labelling function  $\mathcal{L}$  for the source language `Imp` and a program  $P$  in the `Imp` language, we define an annotation for the source program as follows:*

$$An_{\text{Imp}}(P) = \mathcal{I}(\mathcal{L}(P)) .$$

**Proposition 8.** *If  $P$  is a program and  $\mathcal{C}'(\mathcal{C}(\mathcal{L}(P)))$  is a sound (sound and precise) labelling then  $(An_{\text{Imp}}(P), s[c/\text{cost}]) \Downarrow s'[c + \delta/\text{cost}]$  and  $m \Vdash \epsilon, s[c/\text{cost}]$  entails that  $(\mathcal{C}'(\mathcal{C}(P)), m) \Downarrow m', m' \Vdash \epsilon, s'[c/\text{cost}]$  and the cost of the execution is bounded by (is exactly)  $\delta$ .*

### 3 A C compiler producing cost annotations

We now consider an untrusted C compiler prototype in `ocaml` in order to experiment with the scalability of our approach. Its architecture is described below:



The most notable difference with `CompCert` [10] is that we target the Intel 8051 [9] and Mips assembly languages (rather than `PowerPc`). The compilation from C to Clight relies on the CIL front-end [12]. The one from Clight to RTL has been programmed from scratch and it is partly based on the Coq definitions available in the `CompCert` compiler. Finally, the back-end from RTL to Mips is based on a compiler developed in `ocaml` for pedagogical purposes<sup>7</sup>; we extended this back-end to target the Intel 8051. The main optimizations the back-end performs are liveness analysis and register allocation, and graph compression. We ran some benchmarks to ensure that our prototype implementation is realistic. The results are given in appendix B.9.

This section informally describes the labelled extensions of the languages in the compilation chain (see Appendix B for details), the way the labels are propagated by the compilation functions, and the (sound and precise) labelling

<sup>7</sup> <http://www.enseignement.polytechnique.fr/informatique/INF564/>

of the source code. A related experiment concerning a higher-order functional language of the ML family is described in [2].

### 3.1 Labelled languages

Both the `Clight` and `Cminor` languages are extended in the same way by labelling both statements and expressions (by comparison, in the toy language `Imp` we just used labelled statements). The labelling of expressions aims to capture precisely their execution cost. Indeed, `Clight` and `Cminor` include expressions such as  $a_1 ? a_2 ; a_3$  whose evaluation cost depends on the boolean value  $a_1$ . As both languages are extended in the same way, the extended compilation does nothing more than sending `Clight` labelled statements and expressions to those of `Cminor`.

The labelled versions of `RTLAbs` and the languages in the back-end simply consist in adding a new instruction whose semantics is to emit a label without modifying the state. For the CFG based languages (`RTLAbs` to `LTL`), this new instruction is `emit label → node`. For `LIN`, `Mips` and `8051`, it is `emit label`. The translation of these label instructions is immediate.

### 3.2 Labelling of the source language

As for the toy compiler, the goals of a labelling are soundness, precision, and possibly economy. Our labelling for `Clight` resembles that of `Imp` for their common instructions (e.g. loops). We only consider the instructions of `Clight` that are not present in `Imp`<sup>8</sup>.

*Ternary expressions* They may introduce a branching in the control flow. We achieve precision by associating a label with each branch.

*Program Labels and Gotos* Program labels and `gotos` are intraprocedural. Their only effect on the control flow is to potentially introduce an unguarded loop. This loop must contain at least one cost label in order to satisfy the soundness condition, which we ensure by adding a cost label right after each program label.

*Function calls* In the general case, the address of the callee cannot be inferred statically. But in the compiled assembly code, we know for a fact that the callee ends with a return statement that transfers the control back to the instruction following the function call in the caller. As a result, we treat function calls according to the following invariants: (1) the instructions of a function are covered by the labels inside this function, (2) we assume a function call always returns and runs the instruction following the call. Invariant (1) entails in particular that each function must contain at least one label. Invariant (2) is of course an over-approximation of the program behavior as a function might fail to return because of an error or an infinite loop. In this case, the proposed labelling remains correct: it just assumes that the instructions following the function call will be executed, and takes their cost into consideration. The final computed cost is still an over-approximation of the actual cost.

<sup>8</sup> We do not consider expressions with side-effects because they are eliminated by CIL.

## 4 A tool for reasoning on cost annotations

Frama – C is a set of analysers for C programs with a specification language called ACSL. New analyses can be dynamically added through a plug-in system. For instance, the *Jessie* plug-in allows deductive verification of C programs with respect to their specification in ACSL, with various provers as back-end tools.

We developed the *Cost* plug-in for the Frama – C platform as a proof of concept of an automatic environment exploiting the cost annotations produced by the *CerCo* compiler. It consists of an *ocaml* program of 5 *Kloc* which in first approximation takes the following actions: (1) it receives as input a C program, (2) it applies the *CerCo* compiler to produce a related C program with cost annotations, (3) it applies some heuristics to produce a tentative bound on the cost of executing the C functions of the program as a function of the value of their parameters, (4) the user can then call the *Jessie* tool to discharge the related proof obligations. In the following we elaborate on the soundness of the framework, the algorithms underlying the plug-in, and the experiments we performed with the *Cost* tool.

### 4.1 Soundness

The soundness of the whole framework depends on the cost annotations added by the *CerCo* compiler, the synthetic costs produced by the *Cost* plug-in, the verification conditions (VCs) generated by *Jessie*, and the external provers discharging the VCs. The synthetic costs being in ACSL format, *Jessie* can be used to verify them. Thus, even if the added synthetic costs are incorrect (relatively to the cost annotations), the process in its globality is still correct: indeed, *Jessie* will not validate incorrect costs and no conclusion can be made about the WCET of the program in this case. In other terms, the soundness does not really depend on the action of the *Cost* plug-in, which can in principle produce *any* synthetic cost. However, in order to be able to actually prove a WCET of a C function, we need to add correct annotations in a way that *Jessie* and subsequent automatic provers have enough information to deduce their validity. In practice this is not straightforward even for very simple programs composed of branching and assignments (no loops and no recursion) because a fine analysis of the VCs associated with branching may lead to a complexity blow up.

### 4.2 Inner workings

The cost annotations added by the *CerCo* compiler take the form of C instructions that update by a constant a fresh global variable called the *cost variable*. Synthesizing a WCET of a C function thus consists in statically resolving an upper bound of the difference between the value of the cost variable before and after the execution of the function, i.e. find in the function the instructions that update the cost variable and establish the number of times they are passed through during the flow of execution. The plug-in proceeds as follows.

- Each function is independently processed and is associated a WCET that may depend on the cost of the other functions. This is done with a mix between abstract interpretation [5] and syntactic recognition of specific loops.
- As result of the previous step, a system of inequations is built and its solution is attempted by an iterative process. At each iteration, one replaces in all the inequations the references to the cost of a function by its associated cost if it is independent of the other functions. This step is repeated till a fixpoint is reached.
- ACSL annotations are added to the program according to the result of the above fixpoint. The two previous steps may fail in finding a concrete WCET for some functions, because of imprecision inherent in abstract interpretation, and because of recursive definitions in the source program not solved by the fixpoint. At each program point that requires an annotation (function definitions and loops), annotations are added if a solution was found for the program point.
- The most difficult instructions to handle are loops. We consider loops for which we can syntactically find a counter (its initial, increment and last values are domain dependent). Other loops are associated an undefined cost ( $\top$ ). When encountering a loop, the analysis first sets the cost of its entry point to 0. The cost inside the loop is thus relative to the loop. Then, for each exit point, we fetch the value of the cost at that point and multiply it by an upper bound of the number of iterations (obtained through arithmetic over the initial, increment and last values of the counter); this results in an upper bound of the cost of the whole loop, which is sent to the successors of the considered exit point.

Figure 2 shows the action of the *Cost* plug-in on a C program. The most notable differences are the added so-called *cost variable*, some associated update (increment) instructions inside the code, and an **ensures** clause that specifies the WCET of the `is_sorted` function with respect to the cost variable. One can notice that this WCET depends on the inputs of the function. Running *Jessie* on the annotated and specified program generates VCs that are all proved by the automatic prover *AltErgo*<sup>9</sup>.

### 4.3 Experiments

The *Cost* plug-in has been developed in order to validate *CerCo*'s framework for modular WCET analysis. The plug-in allows (semi-)automatic generation and certification of WCET for C programs. Also, we designed a wrapper for supporting *Lustre* files. Indeed, *Lustre* is a data-flow language to program synchronous systems and the language comes with a compiler to C. The C function produced by the compiler implements the *step function* of the synchronous system and computing the WCET of the function amounts to obtain a bound on the reaction time of the system.

<sup>9</sup> *AltErgo* prover. <http://ergo.lri.fr/>

```

int is_sorted (int *tab, int size) {
  int i, res = 1;
  for (i = 0 ; i < size-1 ; i++) if (tab[i] > tab[i+1]) res = 0;
  return res; }

```

(a) The initial C source code.

```

int _cost = 0;

/*@ ensures (_cost ≤ \old(_cost)+(101+(0<size-1?(size-1)*195:0))); */
int is_sorted (int *tab, int size) {
  int i, res = 1, _cost_tmp0;
  _cost += 97; _cost_tmp0 = _cost;
  /*@ loop invariant (0 < size-1) ⇒ (i ≤ size-1);
   @ loop invariant (0 ≥ size-1) ⇒ (i ≡ 0);
   @ loop invariant (_cost ≤ _cost_tmp0+i*195);
   @ loop variant (size-1)-i; */
  for (i = 0; i < size-1; i++) {
    _cost += 91;
    if (tab[i] > tab[i+1]) { _cost += 104; res = 0; }
    else _cost += 84; }
  _cost += 4; return res; }

```

(b) The annotated source code generated by Cost.

Figure2: An example of the Cost plug-in action.

We tested the Cost plug-in and the Lustre wrapper on the C programs generated by the Lustre compiler. We also tested it on some basic algorithms and cryptographic functions; these examples, unlike those generated by the Lustre compiler include arrays and for-loops. Table 3 provides a list of concrete programs and describes their type, functionality, the number of lines of the source code, and the number of VCs generated. In each case, the Cost plug-in computes a WCET and AltErgo is able to discharge all VCs. Obviously the generation of synthetic costs is an undecidable and open-ended problem. Our experience just shows that there are classes of C programs which are relevant for embedded applications and for which the synthesis and verification tasks can be completely automatized.

*Acknowledgement* The master students Kayvan Memarian and Ronan Saillard contributed both to the Coq proofs and the CerCo compiler in the early stages of their development.

File	Type	Description	LOC	VCs
<code>3-way.c</code>	C	Three way block cipher	144	34
<code>a5.c</code>	C	A5 stream cipher, used in GSM cellular	226	18
<code>array_sum.c</code>	S	Sums the elements of an integer array	15	9
<code>fact.c</code>	S	Factorial function, imperative implementation	12	9
<code>is_sorted.c</code>	S	Sorting verification of an array	8	8
<code>LFSR.c</code>	C	32-bit linear-feedback shift register	47	3
<code>minus.c</code>	L	Two modes button	193	8
<code>mmb.c</code>	C	Modular multiplication-based block cipher	124	6
<code>parity.lus</code>	L	Parity bit of a boolean array	359	12
<code>random.c</code>	C	Random number generator	146	3
S: standard algorithm    C: cryptographic function L: C generated from a Lustre file				

Figure3: Experiments on CerCo and the Cost plug-in.

## References

1. R.M. Amadio, N. Ayache, K. Memarian, R. Saillard, Y. Régis-Gianas. Compiler Design and Intermediate Languages. Deliverable 2.1 of [3].
2. R.M. Amadio, Y. Régis-Gianas. Certifying and reasoning on cost annotations of functional programs. In Proc. FOPARA 2011, Springer LNCS 7177, 2012.
3. Certified complexity (Project description). ICT-2007.8.0 FET Open, Grant 243881.
4. L. Correnson, P. Cuoq, F. Kirchner, V. Prevosto, A. Puccetti, J. Signoles, B. Yakobowski. Frama-C user manual. CEA-LIST, Software Safety Laboratory, Saclay, F-91191. <http://frama-c.com/>
5. P. Cousot, R. Cousot. Abstract Interpretation Frameworks. Journal of Logic and Computation, volume 2, number 4, pages 511–547, Oxford University Press, Oxford, UK, 1992.
6. C. Ferdinand, R. Heckmann, T. Le Sergent, D. Lopes, B. Martin, X. Fornari, and F. Martin. Combining a high-level design tool for safety-critical systems with a tool for WCET analysis of executables. In *Embedded Real Time Software*, 2008.
7. X. Fornari. Understanding how SCADE suite KCG generates safe C code. White paper, Esterel Technologies, 2010.
8. J. Larus. Assemblers, linkers, and the SPIM simulator. Appendix of *Computer Organization and Design: the hw/sw interface*, by Hennessy and Patterson, 2005.
9. MCS 51 Microcontroller Family User’s Manual. Publication number 121517, by Intel Corporation, 1994,
10. X. Leroy. Formal verification of a realistic compiler. *Commun. ACM*, 52(7):107–115, 2009.
11. X. Leroy. Mechanized semantics, with applications to program proof and compiler verification. *Marktoberdorf summer school*, 2009.
12. G. Necula, S. McPeak, S.P. Rahul, and W. Weimer. CIL: Intermediate Language and Tools for Analysis and Transformation of C Programs. In *Proceedings of Conference on Compiler Construction*, Springer LNCS 2304:213–228, 2002.
13. R. Wilhelm et al. The worst-case execution-time problem - overview of methods and survey of tools. *ACM Trans. Embedded Comput. Syst.*, 7(3), 2008.



## A A toy compiler

We formalize the toy compiler introduced in section 2 and sketch the related proofs. The related, browsable formal development in COQ is available at <http://www.pps.jussieu.fr/~yrg/cerco>.

### A.1 Imp: language and semantics

The syntax of the **Imp** language is described below. This is a rather standard imperative language with **while** loops and **if-then-else**.

$id ::= x \mid y \mid \dots$	(identifiers)
$n ::= 0 \mid -1 \mid +1 \mid \dots$	(integers)
$v ::= n \mid \text{true} \mid \text{false}$	(values)
$e ::= id \mid n \mid e + e$	(numerical expressions)
$b ::= e < e$	(boolean conditions)
$S ::= \text{skip} \mid id := e \mid S; S \mid \text{if } b \text{ then } S \text{ else } S \mid \text{while } b \text{ do } S$	(commands)
$P ::= \text{prog } S$	(programs)

Let  $s$  be a total function from identifiers to integers representing the **state**. If  $s$  is a state,  $x$  an identifier, and  $n$  an integer then  $s[n/x]$  is the ‘updated’ state such that  $s[n/x](x) = n$  and  $s[n/x](y) = s(y)$  if  $x \neq y$ . The *big-step* operational semantics of **Imp** expressions and boolean conditions is defined as follows:

$$\frac{}{(v, s) \Downarrow v} \quad \frac{}{(x, s) \Downarrow s(x)} \quad \frac{(e, s) \Downarrow v \quad (e', s) \Downarrow v'}{(e + e', s) \Downarrow (v + \mathbf{z} \ v')} \quad \frac{(e, s) \Downarrow v \quad (e', s) \Downarrow v'}{(e < e', s) \Downarrow (v < \mathbf{z} \ v')}$$

A *continuation*  $K$  is a list of commands which terminates with a special symbol **halt**:  $K ::= \text{halt} \mid S \cdot K$ . Table 2 defines a small-step semantics of **Imp** commands whose basic judgement has the shape:  $(S, K, s) \rightarrow (S', K', s')$ . We define the semantics of a program **prog**  $S$  as the semantics of the command  $S$  with continuation **halt**. We derive a big step semantics from the small step one as follows:  $(S, s) \Downarrow s'$  if  $(S, \text{halt}, s) \rightarrow \dots \rightarrow (\text{skip}, \text{halt}, s')$ .

### A.2 Vm: language and semantics

Following [11], we define a virtual machine **Vm** and its programming language. The machine includes the following elements: (1) a fixed code  $C$  (a possibly empty sequence of instructions), (2) a program counter  $pc$ , (3) a store  $s$  (as for the source program), (4) a stack of integers  $\sigma$ .

Given a sequence  $C$ , we denote with  $|C|$  its length and with  $C[i]$  its  $i^{\text{th}}$  element (the leftmost element being the  $0^{\text{th}}$  element). The operational semantics of the instructions is formalized by rules of the shape  $C \vdash (i, \sigma, s) \rightarrow (j, \sigma', s')$  and it is fully described in table 3. Notice that **Imp** and **Vm** semantics share the same notion of store. We write, *e.g.*,  $n \cdot \sigma$  to stress that the top element of the stack

$(x := e, K, s)$	$\rightarrow (\text{skip}, K, s[v/x])$	if $(e, s) \Downarrow v$
$(S; S', K, s)$	$\rightarrow (S, S' \cdot K, s)$	
$(\text{if } b \text{ then } S \text{ else } S', K, s)$	$\rightarrow \begin{cases} (S, K, s) & \text{if } (b, s) \Downarrow \text{true} \\ (S', K, s) & \text{if } (b, s) \Downarrow \text{false} \end{cases}$	
$(\text{while } b \text{ do } S, K, s)$	$\rightarrow \begin{cases} (S, (\text{while } b \text{ do } S) \cdot K, s) & \text{if } (b, s) \Downarrow \text{true} \\ (\text{skip}, K, s) & \text{if } (b, s) \Downarrow \text{false} \end{cases}$	
$(\text{skip}, S \cdot K, s)$	$\rightarrow (S, K, s)$	

Table2: Small-step operational semantics of **Imp** commands

Rule	$C[i] =$
$C \vdash (i, \sigma, s) \rightarrow (i + 1, n \cdot \sigma, s)$	<b>cnst</b> ( $n$ )
$C \vdash (i, \sigma, s) \rightarrow (i + 1, s(x) \cdot \sigma, s)$	<b>var</b> ( $x$ )
$C \vdash (i, n \cdot \sigma, s) \rightarrow (i + 1, \sigma, s[n/x])$	<b>setvar</b> ( $x$ )
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + 1, (n +_{\mathbf{Z}} n') \cdot \sigma, s)$	<b>add</b>
$C \vdash (i, \sigma, s) \rightarrow (i + k + 1, \sigma, s)$	<b>branch</b> ( $k$ )
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + 1, \sigma, s)$	<b>bge</b> ( $k$ ) and $n <_{\mathbf{Z}} n'$
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + k + 1, \sigma, s)$	<b>bge</b> ( $k$ ) and $n \geq_{\mathbf{Z}} n'$

Table3: Operational semantics **Vm** programs

exists and is  $n$ . We will also write  $(C, s) \Downarrow s'$  if  $C \vdash (0, \epsilon, s) \xrightarrow{*} (i, \epsilon, s')$  and  $C[i] = \text{halt}$ .

Code coming from the compilation of **Imp** programs has specific properties that are used in the following compilation step when values on the stack are allocated either in registers or in main memory. In particular, it turns out that for every instruction of the compiled code it is possible to predict statically the *height of the stack* whenever the instruction is executed. We now proceed to define a simple notion of *well-formed* code and show that it enjoys this property. In the following section, we will define the compilation function from **Imp** to **Vm** and show that it produces well-formed code.

**Definition 5.** *We say that a sequence of instructions  $C$  is well formed if there is a function  $h : \{0, \dots, |C|\} \rightarrow \mathbf{N}$  which satisfies the conditions listed in table 4 for  $0 \leq i \leq |C| - 1$ . In this case we write  $C : h$ .*

The conditions defining the predicate  $C : h$  are strong enough to entail that  $h$  correctly predicts the stack height and to guarantee the uniqueness of  $h$  up to the initial condition.

**Proposition 9.** (1) *If  $C : h$ ,  $C \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma', s')$ , and  $h(i) = |\sigma|$  then  $h(j) = |\sigma'|$ . (2) *If  $C : h$ ,  $C : h'$  and  $h(0) = h'(0)$  then  $h = h'$ .**

$C[i] =$	Conditions for $C : h$
$\text{cns}(n)$ or $\text{var}(x)$	$h(i+1) = h(i) + 1$
$\text{add}$	$h(i) \geq 2, \quad h(i+1) = h(i) - 1$
$\text{setvar}(x)$	$h(i) = 1, \quad h(i+1) = 0$
$\text{branch}(k)$	$0 \leq i+k+1 \leq  C , \quad h(i) = h(i+1) = h(i+k+1) = 0$
$\text{bge}(k)$	$0 \leq i+k+1 \leq  C , \quad h(i) = 2, \quad h(i+1) = h(i+k+1) = 0$
$\text{halt}$	$i =  C  - 1, \quad h(i) = h(i+1) = 0$

Table4: Conditions for well-formed code

### A.3 Compilation from **Imp** to **Vm**

In table 5, we define compilation functions  $\mathcal{C}$  from **Imp** to **Vm** which operate on expressions, boolean conditions, statements, and programs. We write  $sz(e)$ ,  $sz(b)$ ,  $sz(S)$  for the number of instructions the compilation function associates with the expression  $e$ , the boolean condition  $b$ , and the statement  $S$ , respectively.

$$\mathcal{C}(x) = \text{var}(x) \quad \mathcal{C}(n) = \text{cns}(n) \quad \mathcal{C}(e + e') = \mathcal{C}(e) \cdot \mathcal{C}(e') \cdot \text{add}$$

$$\mathcal{C}(e < e', k) = \mathcal{C}(e') \cdot \mathcal{C}(e) \cdot \text{bge}(k)$$

$$\mathcal{C}(x := e) = \mathcal{C}(e) \cdot \text{setvar}(x) \quad \mathcal{C}(S; S') = \mathcal{C}(S) \cdot \mathcal{C}(S')$$

$$\mathcal{C}(\text{if } b \text{ then } S \text{ else } S') = \mathcal{C}(b, k) \cdot \mathcal{C}(S) \cdot (\text{branch}(k')) \cdot \mathcal{C}(S') \\ \text{where: } k = sz(S) + 1, \quad k' = sz(S')$$

$$\mathcal{C}(\text{while } b \text{ do } S) = \mathcal{C}(b, k) \cdot \mathcal{C}(S) \cdot \text{branch}(k') \\ \text{where: } k = sz(S) + 1, \quad k' = -(sz(b) + sz(S) + 1)$$

$$\mathcal{C}(\text{prog } S) = \mathcal{C}(S) \cdot \text{halt}$$

Table5: Compilation from **Imp** to **Vm**

We follow [11] for the proof of soundness of the compilation function for expressions and boolean conditions.

**Proposition 10.** *The following properties hold:*

- (1) *If  $(e, s) \Downarrow v$  then  $C \cdot \mathcal{C}(e) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j, v \cdot \sigma, s)$  where  $i = |C|$  and  $j = |C \cdot \mathcal{C}(e)|$ .*
- (2) *If  $(b, s) \Downarrow \text{true}$  then  $C \cdot \mathcal{C}(b, k) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j + k, \sigma, s)$  where  $i = |C|$  and  $j = |C \cdot \mathcal{C}(b, k)|$ .*
- (3) *If  $(b, s) \Downarrow \text{false}$  then  $C \cdot \mathcal{C}(b, k) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s)$  where  $i = |C|$  and  $j = |C \cdot \mathcal{C}(b, k)|$ .*

Next we focus on the compilation of statements. We introduce a ternary relation  $R(C, i, K)$  which relates a Vm code  $C$ , a number  $i \in \{0, \dots, |C| - 1\}$  and a continuation  $K$ . The intuition is that relative to the code  $C$ , the instruction  $i$  can be regarded as having continuation  $K$ . (A formal definition is available in Appendix 11.) We can then state the correctness of the compilation function as follows.

**Proposition 11.** *If  $(S, K, s) \rightarrow (S', K', s')$  and  $R(C, i, S \cdot K)$  then  $C \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s')$  and  $R(C, j, S' \cdot K')$ .*

As announced, we can prove that the result of the compilation is a well-formed code.

**Proposition 12.** *For any program  $P$  there is a unique  $h$  such that  $\mathcal{C}(P) : h$ .*

#### A.4 Mips: language and semantics

We consider a Mips-like machine [8] which includes the following elements: (1) a fixed code  $M$  (a sequence of instructions), (2) a program counter  $pc$ , (3) a finite set of registers including the registers  $A$ ,  $B$ , and  $R_0, \dots, R_{b-1}$ , and (4) an (infinite) main memory which maps locations to integers.

We denote with  $R, R', \dots$  registers, with  $l, l', \dots$  locations and with  $m, m', \dots$  memories which are total functions from registers and locations to (unbounded) integers. We denote with  $M$  a list of instructions. The operational semantics is formalized in table 6 by rules of the shape  $M \vdash (i, m) \rightarrow (j, m')$ , where  $M$  is a list of Mips instructions,  $i, j$  are natural numbers and  $m, m'$  are memories. We write  $(M, m) \Downarrow m'$  if  $M \vdash (0, m) \xrightarrow{*} (j, m')$  and  $M[j] = \text{halt}$ .

Rule	$M[i] =$
$M \vdash (i, m) \rightarrow (i + 1, m[n/R])$	loadi $R, n$
$M \vdash (i, m) \rightarrow (i + 1, m[m(l)/R])$	load $R, l$
$M \vdash (i, m) \rightarrow (i + 1, m[m(R)/l])$	store $R, l$
$M \vdash (i, m) \rightarrow (i + 1, m[m(R') + m(R'')/R])$	add $R, R', R''$
$M \vdash (i, m) \rightarrow (i + k + 1, m)$	branch $k$
$M \vdash (i, m) \rightarrow (i + 1, m)$	bge $R, R', k$ and $m(R) <_{\mathbf{Z}} m(R')$
$M \vdash (i, m) \rightarrow (i + k + 1, m)$	bge $R, R', k$ and $m(R) \geq_{\mathbf{Z}} m(R')$

Table6: Operational semantics Mips programs

#### A.5 Compilation from Vm to Mips

In order to compile Vm programs to Mips programs we make the following hypotheses: (1) for every Vm program variable  $x$  we reserve an address  $l_x$ , (2) for

every natural number  $h \geq b$ , we reserve an address  $l_h$  (the addresses  $l_x, l_h, \dots$  are all distinct), and (3) we store the first  $b$  elements of the stack  $\sigma$  in the registers  $R_0, \dots, R_{b-1}$  and the remaining (if any) at the addresses  $l_b, l_{b+1}, \dots$

We say that the memory  $m$  represents the stack  $\sigma$  and the store  $s$ , and write  $m \Vdash \sigma, s$ , if the following conditions are satisfied: (1)  $s(x) = m(l_x)$ , and (2) if  $0 \leq i < |\sigma|$  then  $\sigma[i] = m(R_i)$  if  $i < b$ , and  $\sigma[i] = m(l_i)$  if  $i \geq b$ .

$C[i] =$	$C'(i, C) =$
$\text{cnst}(n)$	$\begin{cases} (\text{loadi } R_h, n) & \text{if } h = h(i) < b \\ (\text{loadi } A, n) \cdot (\text{store } A, l_h) & \text{otherwise} \end{cases}$
$\text{var}(x)$	$\begin{cases} (\text{load } R_h, l_x) & \text{if } h = h(i) < b \\ (\text{load } A, l_x) \cdot (\text{store } A, l_h) & \text{otherwise} \end{cases}$
$\text{add}$	$\begin{cases} (\text{add } R_{h-2}, R_{h-2}, R_{h-1}) & \text{if } h = h(i) < (b-1) \\ (\text{load } A, l_{h-1}) \cdot (\text{add } R_{h-2}, R_{h-2}, A) & \text{if } h = h(i) = (b-1) \\ (\text{load } A, l_{h-1}) \cdot (\text{load } B, l_{h-2}) & \text{if } h = h(i) > (b-1) \\ (\text{add } A, B, A) \cdot (\text{store } A, l_{h-2}) & \end{cases}$
$\text{setvar}(x)$	$\begin{cases} (\text{store } R_{h-1} l_x) & \text{if } h = h(i) < b \\ (\text{load } A, l_{h-1}) \cdot (\text{store } A, l_x) & \text{if } h = h(i) \geq b \end{cases}$
$\text{branch}(k)$	$(\text{branch } k') \text{ if } k' = p(i + k + 1, C) - p(i + 1, C)$
$\text{bge}(k)$	$\begin{cases} (\text{bge } R_{h-2}, R_{h-1}, k') & \text{if } h = h(i) < (b-1) \\ (\text{load } A, l_{h-1}) \cdot (\text{bge } R_{h-2}, A, k') & \text{if } h = h(i) = (b-1) \\ (\text{load } A, l_{h-2}) \cdot (\text{load } B, l_{h-1}) \cdot (\text{bge } A, B, k') & \text{if } h = h(i) > (b-1), k' = \\ & p(i + k + 1, C) - p(i + 1, C) \end{cases}$
$\text{halt}$	$\text{halt}$

Table7: Compilation from Vm to Mips

The compilation function  $C'$  from Vm to Mips is described in table 7. It operates on a well-formed Vm code  $C$  whose last instruction is **halt**. Hence, by proposition 12(3), there is a unique  $h$  such that  $C : h$ . We denote with  $C'(C)$  the concatenation  $C'(0, C) \dots C'(|C| - 1, C)$ . Given a well formed Vm code  $C$  with  $i < |C|$  we denote with  $p(i, C)$  the position of the first instruction in  $C'(C)$  which corresponds to the compilation of the instruction with position  $i$  in  $C$ . This is defined as<sup>10</sup>  $p(i, C) = \sum_{0 \leq j < i} d(j, C)$ , where the function  $d(i, C)$  is defined as  $d(i, C) = |C'(i, C)|$ . Hence  $d(i, C)$  is the number of Mips instructions associated with the  $i^{\text{th}}$  instruction of the (well-formed)  $C$  code. The functional correctness of the compilation function can then be stated as follows.

**Proposition 13.** *Let  $C : h$  be a well formed code. If  $C \vdash (i, \sigma, s) \rightarrow (j, \sigma', s')$  with  $h(i) = |\sigma|$  and  $m \Vdash \sigma, s$  then  $C'(C) \vdash (p(i, C), m) \xrightarrow{*} (p(j, C), m')$  and  $m' \Vdash \sigma', s'$ .*

<sup>10</sup> There is an obvious circularity in this definition that can be easily eliminated by defining first the function  $d$  following the case analysis in table 7, then the function  $p$ , and finally the function  $C'$  as in table 7.

## A.6 Notation

Let  $\xrightarrow{t}$  be a family of reduction relations where  $t$  ranges over the set of labels and  $\epsilon$ . Then we define:

$$\xRightarrow{t} = \begin{cases} (\xrightarrow{\epsilon})^* & \text{if } t = \epsilon \\ (\xrightarrow{\epsilon})^* \circ \xrightarrow{t} \circ (\xrightarrow{\epsilon})^* & \text{otherwise} \end{cases}$$

where as usual  $R^*$  denote the reflexive and transitive closure of the relation  $R$  and  $\circ$  denotes the composition of relations.

## A.7 Proof of proposition 1

- (1) By induction on the structure of the command  $S$ .
- (2) By iterating the following proposition.

**Proposition 14.** *If  $(S, K, s) \xrightarrow{t} (S', K', s')$  and  $R(C, i, S \cdot K)$  with  $t = \ell$  or  $t = \epsilon$  then  $C \vdash (i, \sigma, s) \xRightarrow{t} (j, \sigma, s')$  and  $R(C, j, S' \cdot K')$ .*

This is an extension of proposition 11 and it is proven in the same way with an additional case for labelled commands.  $\square$

## A.8 Proof of proposition 2

- (1) The compilation of the `Vm` instruction `emit( $\ell$ )` is the `Mips` instruction `(emit  $\ell$ )`.
- (2) By iterating the following proposition.

**Proposition 15.** *Let  $C : h$  be a well formed code. If  $C \vdash (i, \sigma, s) \xrightarrow{t} (j, \sigma', s')$  with  $t = \ell$  or  $t = \epsilon$ ,  $h(i) = |\sigma|$  and  $m \Vdash \sigma, s$  then  $C'(C) \vdash (p(i, C), m) \xRightarrow{t} (p(j, C), m')$  and  $m' \Vdash \sigma', s'$ .*

## A.9 Proof of proposition 3

We extend the instrumentation to the continuations by defining:

$$\mathcal{I}(S \cdot K) = \mathcal{I}(S) \cdot \mathcal{I}(K) \quad \mathcal{I}(\text{halt}) = \text{halt} .$$

Then we examine the possible reductions of a configuration  $(\mathcal{I}(S), \mathcal{I}(K), s[c/\text{cost}])$ .

- If  $S$  is an unlabelled statement such as `while  $b$  do  $S'$`  then  $\mathcal{I}(S) = \text{while } b \text{ do } \mathcal{I}(S')$  and assuming  $(b, s) \Downarrow \text{true}$  the reduction step is:

$$(\mathcal{I}(S), \mathcal{I}(K), s[c/\text{cost}]) \rightarrow (\mathcal{I}(S'), \mathcal{I}(S) \cdot \mathcal{I}(K), s[c/\text{cost}]) .$$

Noticing that  $\mathcal{I}(S) \cdot \mathcal{I}(K) = \mathcal{I}(S \cdot K)$ , this step is matched in the labelled language as follows:

$$(S, K, s[c/\text{cost}]) \rightarrow (S', S \cdot K, s[c/\text{cost}]) .$$

- On the other hand, if  $S = \ell : S'$  is a labelled statement then  $\mathcal{I}(S) = inc(\ell); \mathcal{I}(S')$  and, by a sequence of reductions steps, we have:

$$(\mathcal{I}(S), \mathcal{I}(K), s[c/cost]) \xrightarrow{*} (\mathcal{I}(S'), \mathcal{I}(K), s[c + \kappa(\ell)/cost]) .$$

This step is matched by the labelled reduction:

$$(S, K, s[c/cost]) \xrightarrow{\ell} (S', K, s[c/cost]) .$$

□

#### A.10 Proof of proposition 4

By diagram chasing using propositions 1(1), 2(1), and the definition 1 of labelling. □

#### A.11 Proof of proposition 5

Suppose that:

$$(\mathcal{I}(\mathcal{L}(P)), s[c/cost]) \Downarrow s'[c + \delta/cost] \text{ and } m \Vdash s[c/cost] .$$

Then, by proposition 3, for some  $\lambda$ :

$$(\mathcal{L}(P), s[c/cost]) \Downarrow (s'[c/cost], \lambda) \text{ and } \kappa(\lambda) = \delta .$$

Finally, by propositions 1(2) and 2(2) :

$$(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P))), m) \Downarrow (m', \lambda) \text{ and } m' \Vdash s'[c/cost] .$$

□

#### A.12 Proof of proposition 6

We discuss first the case for sound labelling. If  $\lambda = \ell_1 \cdots \ell_n$  then the computation is the concatenation of simple paths labelled with  $\ell_1, \dots, \ell_n$ . Since  $\kappa(\ell_i)$  bounds the cost of a simple path labelled with  $\ell_i$ , the cost of the overall computation is bounded by  $\kappa(\lambda) = \kappa(\ell_1) + \cdots \kappa(\ell_n)$ . For the sound *and precise* labelling the proof above is repeated, by replacing the word *bounds* by *is exactly* and the words *bounded by* by *exactly*. □

#### A.13 Proof of proposition 7

In both labellings under consideration the root node is labelled. An obvious observation is that only commands of the shape **while**  $b$  **do**  $S$  introduce loops in the compiled code. We notice that both labelling introduce a label in the loop (though at different places). Thus all loops go through a label and the compiled code is always sound.

To show the precision of the second labelling  $\mathcal{L}_p$ , we note the following property.

**Lemma 1.** *A soundly labelled graph is precise if each label occurs at most once in the graph and if the immediate successors of the bge nodes are either halt (no successor) or labelled nodes.*

Indeed, in a such a graph starting from a labelled node we can follow a unique path up to a leaf, another labelled node, or a bge node. In the last case, the hypotheses in the lemma 1 guarantee that the two simple paths one can follow from the bge node have the same length/cost.  $\square$

#### A.14 Proof of proposition 8

By applying consecutively propositions 5 and propositions 6.  $\square$

#### A.15 Proof of proposition 11

Given a Vm code  $C$ , we define an ‘accessibility relation’  $\overset{\mathcal{C}}{\sim}$  as the least binary relation on  $\{0, \dots, |C| - 1\}$  such that:

$$\frac{}{i \overset{\mathcal{C}}{\sim} i} \quad \frac{C[i] = \text{branch}(k) \quad (i + k + 1) \overset{\mathcal{C}}{\sim} j}{i \overset{\mathcal{C}}{\sim} j}$$

We also introduce a ternary relation  $R(C, i, K)$  which relates a Vm code  $C$ , a number  $i \in \{0, \dots, |C| - 1\}$  and a continuation  $K$ . The relation is defined as the least one that satisfies the following conditions.

$$\frac{i \overset{\mathcal{C}}{\sim} j \quad C[j] = \text{halt}}{R(C, i, \text{halt})} \quad \frac{i \overset{\mathcal{C}}{\sim} i' \quad C = C_1 \cdot \mathcal{C}(S) \cdot C_2 \quad i' = |C_1| \quad j = |C_1 \cdot \mathcal{C}(S)| \quad R(C, j, K)}{R(C, i, S \cdot K)} .$$

The following properties are useful.

**Lemma 2.** (1) *The relation  $\overset{\mathcal{C}}{\sim}$  is transitive.*

(2) *If  $i \overset{\mathcal{C}}{\sim} j$  and  $R(C, j, K)$  then  $R(C, i, K)$ .*

The first property can be proven by induction on the definition of  $\overset{\mathcal{C}}{\sim}$  and the second by induction on the structure of  $K$ .

Next we can focus on the proposition. The notation  $C \overset{i}{\vdash} C'$  means that  $i = |C|$ . Suppose that:

$$(S, K, s) \rightarrow (S', K', s') \quad (1) \text{ and } R(C, i, S \cdot K) \quad (2) .$$

From (2), we know that there exist  $i'$  and  $i''$  such that:

$$i \overset{\mathcal{C}}{\sim} i' \quad (3), \quad C = C_1 \overset{i'}{\vdash} \mathcal{C}(S) \overset{i''}{\vdash} C_2 \quad (4), \text{ and } R(C, i'', K) \quad (5)$$

and from (3) it follows that:

$$C \vdash (i, \sigma, s) \xrightarrow{*} (i', \sigma, s) \quad (3') .$$



We are looking for  $j$  such that:

$$C \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s') \quad (6), \text{ and } R(C, j, S' \cdot K') \quad (7) .$$

We proceed by case analysis on  $S$ . We just detail the case of the conditional command as the remaining cases have similar proofs. If  $S = \text{if } e_1 < e_2 \text{ then } S_1 \text{ else } S_2$  then (4) is rewritten as follows:

$$C = C_1 \stackrel{i'}{\cdot} \mathcal{C}(e_1) \cdot \mathcal{C}(e_2) \cdot \text{bge}(k_1) \stackrel{a}{\cdot} \mathcal{C}(S_1) \stackrel{b}{\cdot} \text{branch}(k_2) \stackrel{c}{\cdot} \mathcal{C}(S_2) \stackrel{i''}{\cdot} C_2$$

where  $c = a + k_1$  and  $i'' = c + k_2$ . We distinguish two cases according to the evaluation of the boolean condition. We describe the case  $(e_1 < e_2) \Downarrow \text{true}$ . We set  $j = a$ .

- The instance of (1) is  $(S, K, s) \rightarrow (S_1, K, s)$ .
- The reduction required in (6) takes the form  $C \vdash (i, \sigma, s) \xrightarrow{*} (i', \sigma, s) \xrightarrow{*} (a, \sigma, s')$ , and it follows from (3'), the fact that  $(e_1 < e_2) \Downarrow \text{true}$ , and proposition 10(2).
- Property (7), follows from lemma 2(2), fact (5), and the following proof tree:

$$\frac{j \stackrel{c}{\sim} j \quad \frac{b \stackrel{c}{\sim} i'' \quad R(C, i'', K)}{R(C, b, K)}}{R(C, j, S_1 \cdot K)} .$$

□

## B A C compiler

This section gives an informal overview of the compiler, in particular it highlights the main features of the intermediate languages, the purpose of the compilation steps, and the optimisations. The CerCo compiler can be both downloaded and tested as a web application at the URL <http://www.pps.jussieu.fr/~yrg/cerco>.

### B.1 Clight

Clight is a large subset of the C language that we adopt as the source language of our compiler. It features most of the types and operators of C. It includes pointer arithmetic, pointers to functions, and **struct** and **union** types, as well as all C control structures. The main difference with the C language is that Clight expressions are side-effect free, which means that side-effect operators (`=`, `+=`, `++`, ...) and function calls within expressions are not supported. Given a C program, we rely on the CIL tool [12] to deal with the idiosyncrasy of C concrete syntax and to produce an equivalent program in Clight abstract syntax. We refer to the CompCert project [10] for a formal definition of the Clight language. Here we just recall in figure 4 its syntax which is classically structured in expressions, statements, functions, and whole programs. In order to limit the implementation effort, our current compiler for Clight does *not* cover the operators relating to the floating point type `float`. So, in a nutshell, the fragment of C we have implemented is Clight without floating point.

There is a notable difficulty to compile the C language into 8051 assembly code due to the fact that 8051's machine word are 8bits long and C has 32bits primitive datatypes. To deal with the dissimilarity in that case, we first translate the Clight input program into a Clight input program that only uses 8bits primitive datatypes.

### B.2 Cminor

Cminor is a simple, low-level imperative language, comparable to a stripped-down, typeless variant of C. Again we refer to the CompCert project for its formal definition and we just recall in figure 5 its syntax which as for Clight is structured in expressions, statements, functions, and whole programs.

*Translation of Clight to Cminor* As in Cminor stack operations are made explicit, one has to know which variables are stored in the stack. This information is produced by a static analysis that determines the variables whose address may be 'taken'. Also space is reserved for local arrays and structures. In a second step, the proper compilation is performed: it consists mainly in translating Clight control structures to the basic ones available in Cminor.

Expressions:	$a ::=$ $id$ variable identifier $  n$ integer constant $  \text{sizeof}(\tau)$ size of a type $  op_1 a$ unary arithmetic operation $  a op_2 a$ binary arithmetic operation $  *a$ pointer dereferencing $  a.id$ field access $  \&a$ taking the address of $  (\tau)a$ type cast $  a?a : a$ conditional expression
Statements:	$s ::=$ $\text{skip}$ empty statement $  a = a$ assignment $  a = a(a^*)$ function call $  a(a^*)$ procedure call $  s; s$ sequence $  \text{if } a \text{ then } s \text{ else } s$ conditional $  \text{switch } a \text{ sw}$ multi-way branch $  \text{while } a \text{ do } s$ “while” loop $  \text{do } s \text{ while } a$ “do” loop $  \text{for}(s,a,s) s$ “for” loop $  \text{break}$ exit from current loop $  \text{continue}$ next iteration of the current loop $  \text{return } a^?$ return from current function $  \text{goto } lbl$ branching $  lbl : s$ labelled statement
Switch cases:	$sw ::=$ $\text{default} : s$ default case $  \text{case } n : s; sw$ labelled case
Variable declarations:	$dcl ::= (\tau \ id)^*$ type and name
Functions:	$Fd ::= \tau \ id(dcl)\{dcl; s\}$ internal function $  \text{extern } \tau \ id(dcl)$ external function
Programs:	$P ::= dcl; Fd^*; \text{main} = id$ global variables, functions, entry point

Figure4: Syntax of the Clight language

Signatures:	$sig ::= sig\ int\ (int void)$	arguments and result
Expressions:	$a ::=$ <ul style="list-style-type: none"> <li><math>id</math></li> <li><math>  n</math></li> <li><math>  \text{addrsymbol}(id)</math></li> <li><math>  \text{addrstack}(\delta)</math></li> <li><math>  op_1\ a</math></li> <li><math>  op_2\ a\ a</math></li> <li><math>  \kappa[a]</math></li> <li><math>  a?a : a</math></li> </ul>	<ul style="list-style-type: none"> <li>local variable</li> <li>integer constant</li> <li>address of global symbol</li> <li>address within stack data</li> <li>unary arithmetic operation</li> <li>binary arithmetic operation</li> <li>memory read</li> <li>conditional expression</li> </ul>
Statements:	$s ::=$ <ul style="list-style-type: none"> <li><b>skip</b></li> <li><math>  id = a</math></li> <li><math>  \kappa[a] = a</math></li> <li><math>  id^? = a(a) : sig</math></li> <li><math>  \text{tailcall } a(a) : sig</math></li> <li><math>  \text{return}(a^?)</math></li> <li><math>  s; s</math></li> <li><math>  \text{if } a \text{ then } s \text{ else } s</math></li> <li><math>  \text{loop } s</math></li> <li><math>  \text{block } s</math></li> <li><math>  \text{exit } n</math></li> <li><math>  \text{switch } a\ tbl</math></li> <li><math>  lbl : s</math></li> <li><math>  \text{goto } lbl</math></li> </ul>	<ul style="list-style-type: none"> <li>empty statement</li> <li>assignment</li> <li>memory write</li> <li>function call</li> <li>function tail call</li> <li>function return</li> <li>sequence</li> <li>conditional</li> <li>infinite loop</li> <li>block delimiting <b>exit</b> constructs</li> <li>terminate the <math>(n + 1)^{th}</math> enclosing block</li> <li>multi-way test and exit</li> <li>labelled statement</li> <li>jump to a label</li> </ul>
Switch tables:	$tbl ::=$ <ul style="list-style-type: none"> <li><b>default:exit</b>(<math>n</math>)</li> <li><math>  \text{case } i: \text{exit}(n);tbl</math></li> </ul>	
Functions:	$Fd ::=$ <ul style="list-style-type: none"> <li><b>internal</b> <math>sig\ id\ id\ n\ s</math></li> <li><math>  \text{external } id\ sig</math></li> </ul>	<ul style="list-style-type: none"> <li>internal function: signature, parameters, local variables, stack size and body</li> <li>external function</li> </ul>
Programs:	$P ::= \text{prog } (id = data)^* (id = Fd)^* id$	global variables, functions and entry point

Figure5: Syntax of the Cminor language

### B.3 RTLabs

RTLabs is the last architecture independent language in the compilation process. It is a rather straightforward *abstraction* of the *architecture-dependent* RTL intermediate language available in the CompCert project and it is intended to factorize some work common to the various target assembly languages (e.g. optimizations) and thus to make retargeting of the compiler a simpler matter.

We stress that in RTLabs the structure of Cminor expressions is lost and that this may have a negative impact on the following instruction selection step. Still, the subtleties of instruction selection seem rather orthogonal to our goals and we deem the possibility of retargeting easily the compiler more important than the efficiency of the generated code.

*Syntax.* In RTLabs, programs are represented as *control flow graphs* (CFGs for short). We associate with the nodes of the graphs instructions reflecting the Cminor commands. As usual, commands that change the control flow of the program (e.g. loops, conditionals) are translated by inserting suitable branching instructions in the CFG. The syntax of the language is depicted in table 8. Local variables are now represented by *pseudo registers* that are available in unbounded number. The grammar rule *op* that is not detailed in table 8 defines usual arithmetic and boolean operations (+, xor, ≤, etc.) as well as constants and conversions between sized integers.

*Translation of Cminor to RTLabs.* Translating Cminor programs to RTLabs programs mainly consists in transforming Cminor commands in CFGs. Most commands are sequential and have a rather straightforward linear translation. A conditional is translated in a branch instruction; a loop is translated using a back edge in the CFG.

### B.4 RTL

As in RTLabs, the structure of RTL programs is based on CFGs. RTL is the first architecture-dependant intermediate language of our compiler which, in its current version, targets the Mips and 8051 assembly languages.

*Syntax.* RTL is very close to RTLabs. It is based on CFGs and explicits the Mips or the 8051 instructions corresponding to the RTLabs instructions. Type information disappears: everything is represented using machine words. Moreover, each global of the program is associated to an offset. The syntax of the language can be found in table 9. The grammar rules *unop*, *binop*, *uncon*, and *bincon*, respectively, represent the sets of unary operations, binary operations, unary conditions and binary conditions of the target assembly language.

*Translation of RTLabs to RTL.* This translation is mostly straightforward. A RTLabs instruction is often directly translated to a corresponding assembly instruction. There are a few exceptions: some RTLabs instructions are expanded

$$\begin{aligned}
& \text{return\_type} ::= \text{int} \mid \text{void} & \text{signature} ::= (\text{int} \rightarrow)^* \text{return\_type} \\
\\
& \text{memq} ::= \text{int8s} \mid \text{int8u} \mid \text{int16s} \mid \text{int16u} \mid \text{int32} & \text{fun\_ref} ::= \text{fun\_name} \mid \text{psd\_reg} \\
\\
& \text{instruction} ::= \mid \text{skip} \rightarrow \text{node} & \text{(no instruction)} \\
& \quad \mid \text{psd\_reg} := \text{op}(\text{psd\_reg}^*) \rightarrow \text{node} & \text{(operation)} \\
& \quad \mid \text{psd\_reg} := \&\text{var\_name} \rightarrow \text{node} & \text{(address of a global)} \\
& \quad \mid \text{psd\_reg} := \&\text{locals}[n] \rightarrow \text{node} & \text{(address of a local)} \\
& \quad \mid \text{psd\_reg} := \text{fun\_name} \rightarrow \text{node} & \text{(address of a function)} \\
& \quad \mid \text{psd\_reg} := \text{memq}(\text{psd\_reg}[\text{psd\_reg}]) \rightarrow \text{node} & \text{(memory load)} \\
& \quad \mid \text{memq}(\text{psd\_reg}[\text{psd\_reg}]) := \text{psd\_reg} \rightarrow \text{node} & \text{(memory store)} \\
& \quad \mid \text{psd\_reg} := \text{fun\_ref}(\text{psd\_reg}^*) : \text{signature} \rightarrow \text{node} & \text{(function call)} \\
& \quad \mid \text{fun\_ref}(\text{psd\_reg}^*) : \text{signature} & \text{(function tail call)} \\
& \quad \mid \text{test } \text{op}(\text{psd\_reg}^*) \rightarrow \text{node}, \text{node} & \text{(branch)} \\
& \quad \mid \text{return } \text{psd\_reg}? & \text{(return)} \\
\\
& \text{fun\_def} ::= \text{fun\_name}(\text{psd\_reg}^*) : \text{signature} \\
& \quad \text{result} : \text{psd\_reg}? \\
& \quad \text{locals} : \text{psd\_reg}^* \\
& \quad \text{stack} : n \\
& \quad \text{entry} : \text{node} \\
& \quad \text{exit} : \text{node} \\
& \quad (\text{node} : \text{instruction})^* \\
\\
& \text{init\_datum} ::= \text{reserve}(n) \mid \text{int8}(n) \mid \text{int16}(n) \mid \text{int32}(n) & \text{init\_data} ::= \text{init\_datum}^+ \\
\\
& \text{global\_decl} ::= \text{var } \text{var\_name}\{\text{init\_data}\} & \text{fun\_decl} ::= \text{extern } \text{fun\_name}(\text{signature}) \mid \text{fun\_def} \\
\\
& \text{program} ::= \text{global\_decl}^* \\
& \quad \text{fun\_decl}^*
\end{aligned}$$

Table8: Syntax of the RTLabs language

$$\begin{array}{l}
size ::= \text{Byte} \mid \text{HalfWord} \mid \text{Word} \qquad fun\_ref ::= fun\_name \mid psd\_reg \\
\\
instruction ::= \begin{array}{l}
\mid skip \rightarrow node \qquad \qquad \qquad (no \text{ instruction}) \\
\mid psd\_reg := n \rightarrow node \qquad \qquad \qquad (constant) \\
\mid psd\_reg := unop(psd\_reg) \rightarrow node \qquad \qquad \qquad (unary \text{ operation}) \\
\mid psd\_reg := binop(psd\_reg, psd\_reg) \rightarrow node \qquad \qquad \qquad (binary \text{ operation}) \\
\mid psd\_reg := \&globals[n] \rightarrow node \qquad \qquad \qquad (address \text{ of a global}) \\
\mid psd\_reg := \&locals[n] \rightarrow node \qquad \qquad \qquad (address \text{ of a local}) \\
\mid psd\_reg := fun\_name \rightarrow node \qquad \qquad \qquad (address \text{ of a function}) \\
\mid psd\_reg := size(psd\_reg[n]) \rightarrow node \qquad \qquad \qquad (memory \text{ load}) \\
\mid size(psd\_reg[n]) := psd\_reg \rightarrow node \qquad \qquad \qquad (memory \text{ store}) \\
\mid psd\_reg := fun\_ref(psd\_reg^*) \rightarrow node \qquad \qquad \qquad (function \text{ call}) \\
\mid fun\_ref(psd\_reg^*) \qquad \qquad \qquad (function \text{ tail call}) \\
\mid test \text{ uncon}(psd\_reg) \rightarrow node, node \qquad \qquad \qquad (branch \text{ unary condition}) \\
\mid test \text{ bincon}(psd\_reg, psd\_reg) \rightarrow node, node \qquad \qquad \qquad (branch \text{ binary condition}) \\
\mid return \text{ psd\_reg?} \qquad \qquad \qquad (return)
\end{array} \\
\\
fun\_def ::= fun\_name(psd\_reg^*) \qquad \qquad \qquad program ::= \begin{array}{l}
globals : n \\
result : psd\_reg? \\
locals : psd\_reg^* \\
stack : n \\
entry : node \\
exit : node \\
(node : instruction)^*
\end{array} \qquad \qquad \qquad fun\_def^*
\end{array}$$

Table9: Syntax of the RTL language

in two or more assembly instructions. When the translation of a RTLabs instruction requires more than a few simple assembly instruction, it is translated into a call to a function defined in the preamble of the compilation result.

## B.5 ERTL

As in RTL, the structure of ERTL programs is based on CFGs. ERTL explicits the calling conventions of the Mips assembly language. In the back-end for 8051, we defined our own calling convention since there were none.

*Syntax.* The syntax of the language is given in table 10. The main difference between RTL and ERTL is the use of hardware registers. Parameters are passed in specific hardware registers; if there are too many parameters, the remaining are stored in the stack. Other conventionally specific hardware registers are used: a register that holds the result of a function, a register that holds the base address of the globals, a register that holds the address of the top of the stack, and some registers that need to be saved when entering a function and whose values are restored when leaving a function. Following these conventions, function calls do not list their parameters anymore; they only mention their number. Two new instructions appear to allocate and deallocate on the stack some space needed by a function to execute. Along with these two instructions come two instructions to fetch or assign a value in the parameter sections of the stack; these instructions cannot yet be translated using regular load and store instructions because we do not know the final size of the stack area of each function. At last, the return instruction has a boolean argument that tells whether the result of the function may later be used or not (this is exploited for optimizations).

*Translation of RTL to ERTL.* The work consists in expliciting the conventions previously mentioned. These conventions appear when entering, calling and leaving a function, and when referencing a global variable or the address of a local variable.

*Optimizations.* A *liveness analysis* is performed on ERTL to replace unused instructions by a `skip`. An instruction is tagged as unused when it performs an assignment on a register that will not be read afterwards. Also, the result of the liveness analysis is exploited by a *register allocation* algorithm whose result is to efficiently associate a physical location (a hardware register or an address in the stack) to each pseudo register of the program.

## B.6 LTL

As in ERTL, the structure of LTL programs is based on CFGs. Pseudo registers are not used anymore; instead, they are replaced by physical locations (a hardware register or an address in the stack).



$$size ::= \text{Byte} \mid \text{HalfWord} \mid \text{Word} \qquad fun\_ref ::= fun\_name \mid psd\_reg$$

$instruction ::=$	$skip \rightarrow node$	(no instruction)
	$NewFrame \rightarrow node$	(frame creation)
	$DelFrame \rightarrow node$	(frame deletion)
	$psd\_reg := stack[slot, n] \rightarrow node$	(stack load)
	$stack[slot, n] := psd\_reg \rightarrow node$	(stack store)
	$hwd\_reg := psd\_reg \rightarrow node$	(pseudo to hardware)
	$psd\_reg := hwd\_reg \rightarrow node$	(hardware to pseudo)
	$psd\_reg := n \rightarrow node$	(constant)
	$psd\_reg := unop(psd\_reg) \rightarrow node$	(unary operation)
	$psd\_reg := binop(psd\_reg, psd\_reg) \rightarrow node$	(binary operation)
	$psd\_reg := fun\_name \rightarrow node$	(address of a function)
	$psd\_reg := size(psd\_reg[n]) \rightarrow node$	(memory load)
	$size(psd\_reg[n]) := psd\_reg \rightarrow node$	(memory store)
	$fun\_ref(n) \rightarrow node$	(function call)
	$fun\_ref(n)$	(function tail call)
	$test\ uncon(psd\_reg) \rightarrow node, node$	(branch unary condition)
	$test\ bincon(psd\_reg, psd\_reg) \rightarrow node, node$	(branch binary condition)
	$return\ b$	(return)

$$\begin{array}{ll} fun\_def ::= fun\_name(n) & program ::= globals : n \\ & locals : psd\_reg^* \\ & stack : n \\ & entry : node \\ & (node : instruction)^* \end{array}$$

Table10: Syntax of the ERTL language

*Syntax.* Except for a few exceptions, the instructions of the language are those of ERTL with hardware registers replacing pseudo registers. Calling and returning conventions were explicited in ERTL; thus, function calls and returns do not need parameters in LTL. The syntax is defined in table 11.

$size ::= \text{Byte} \mid \text{HalfWord} \mid \text{Word}$		$fun\_ref ::= fun\_name \mid hdw\_reg$	
$instruction ::=$		$\text{skip} \rightarrow node$	(no instruction)
		$\text{NewFrame} \rightarrow node$	(frame creation)
		$\text{DelFrame} \rightarrow node$	(frame deletion)
		$hdw\_reg := n \rightarrow node$	(constant)
		$hdw\_reg := unop(hdw\_reg) \rightarrow node$	(unary operation)
		$hdw\_reg := binop(hdw\_reg, hdw\_reg) \rightarrow node$	(binary operation)
		$hdw\_reg := fun\_name \rightarrow node$	(address of a function)
		$hdw\_reg := size(hdw\_reg[n]) \rightarrow node$	(memory load)
		$size(hdw\_reg[n]) := hdw\_reg \rightarrow node$	(memory store)
		$fun\_ref() \rightarrow node$	(function call)
		$fun\_ref()$	(function tail call)
		$\text{test } uncon(hdw\_reg) \rightarrow node, node$	(branch unary condition)
		$\text{test } bincon(hdw\_reg, hdw\_reg) \rightarrow node, node$	(branch binary condition)
		$\text{return}$	(return)
$fun\_def ::= fun\_name(n)$		$program ::= \text{globals} : n$	
		$\text{locals} : n$	
		$\text{stack} : n$	
		$\text{entry} : node$	
		$(node : instruction)^*$	
		$fun\_def^*$	

Table11: Syntax of the LTL language

*Translation of ERTL to LTL.* The translation relies on the results of the liveness analysis and of the register allocation. Unused instructions are eliminated and each pseudo register is replaced by a physical location. In LTL, the size of the stack frame of a function is known; instructions intended to load or store values in the stack are translated using regular load and store instructions.

*Optimizations.* A *graph compression* algorithm removes empty instructions generated by previous compilation passes and by the liveness analysis.

## B.7 LIN

In LIN, the structure of a program is no longer based on CFGs. Every function is represented as a sequence of instructions.

*Syntax.* The instructions of LIN are very close to those of LTL. *Program labels*, *gotos* and branch instructions handle the changes in the control flow. The syntax of LIN programs is shown in table 12.

$size ::= \text{Byte} \mid \text{HalfWord} \mid \text{Word}$		$fun\_ref ::= fun\_name \mid hdw\_reg$	
$instruction ::=$		$\mid \text{NewFrame}$	(frame creation)
		$\mid \text{DelFrame}$	(frame deletion)
		$\mid hdw\_reg := n$	(constant)
		$\mid hdw\_reg := unop(hdw\_reg)$	(unary operation)
		$\mid hdw\_reg := binop(hdw\_reg, hdw\_reg)$	(binary operation)
		$\mid hdw\_reg := fun\_name$	(address of a function)
		$\mid hdw\_reg := size(hdw\_reg[n])$	(memory load)
		$\mid size(hdw\_reg[n]) := hdw\_reg$	(memory store)
		$\mid \text{call } fun\_ref$	(function call)
		$\mid \text{tailcall } fun\_ref$	(function tail call)
		$\mid uncon(hdw\_reg) \rightarrow node$	(branch unary condition)
		$\mid bincon(hdw\_reg, hdw\_reg) \rightarrow node$	(branch binary condition)
		$\mid asm\_label :$	(assembly label)
		$\mid \text{goto } mips\_label$	(goto)
		$\mid \text{return}$	(return)
$fun\_def ::= fun\_name(n)$		$program ::= \text{globals} : n$	
		$\quad \text{locals} : n$	
		$\quad instruction^*$	
		$\quad fun\_def^*$	

Table12: Syntax of the LIN language

*Translation of LTL to LIN.* This translation amounts to transform in an efficient way the graph structure of functions into a linear structure of sequential instructions.

## B.8 Assembly

We refer to [9] for a description of the 8051 assembly language. In order to compute tight bounds, the CerCo compiler actually generates 8051 binaries. For instance, the number of memory words associated with branching instructions can only be determined at this level (the number depends on the length of the offset).

In the following we describe the simpler Mips assembly language. A program in Mips is a sequence of instructions. The Mips code produced by the compilation of a Clight program starts with a preamble in which some useful and

non-primitive functions are predefined (e.g. conversion from 8 bits unsigned integers to 32 bits integers). The subset of the Mips assembly language that the compilation produces is defined in table 13.

$load ::= lb \mid lhw \mid lw$        $store ::= sb \mid shw \mid sw$        $fun\_ref ::= fun\_name \mid hdw\_reg$

$instruction ::=$	<b>nop</b>	(empty instruction)
	<b>li</b> $hdw\_reg, n$	(constant)
	<b>unop</b> $hdw\_reg, hdw\_reg$	(unary operation)
	<b>binop</b> $hdw\_reg, hdw\_reg, hdw\_reg$	(binary operation)
	<b>la</b> $hdw\_reg, fun\_name$	(address of a function)
	<b>load</b> $hdw\_reg, n(hdw\_reg)$	(memory load)
	<b>store</b> $hdw\_reg, n(hdw\_reg)$	(memory store)
	<b>call</b> $fun\_ref$	(function call)
	<b>uncon</b> $hdw\_reg, node$	(branch unary condition)
	<b>bincon</b> $hdw\_reg, hdw\_reg, node$	(branch binary condition)
	<b>mips_label :</b>	(Mips label)
	<b>j</b> $mips\_label$	(goto)
	<b>return</b>	(return)

$program ::=$	<b>globals :</b> $n$
	<b>entry :</b> $mips\_label^*$
	$instruction^*$

Table13: Syntax of the Mips language

*Translation of LIN to Mips.* This final translation is simple enough. Stack allocation and deallocation are explicited and the function definitions are sequentialized.

## B.9 Benchmarks

To ensure that our prototype compiler is realistic, we performed some preliminary benchmarks on a 183MHz MIPS 4KEc processor, running a linux based distribution. We compared the wall clock execution time of several simple C programs compiled with our compiler against the ones produced by Gcc set up with optimization levels 0 and 1. As shown by Figure 6, our prototype compiler produces executable programs that are on average faster than Gcc's without optimizations.

	gcc -O0	acc gcc -O1	
badsort	55.93	34.51	12.96
fib	76.24	34.28	45.68
mat_det	163.42	156.20	54.76
min	12.21	16.25	3.95
quicksort	27.46	17.95	9.41
search	463.19	623.79	155.38

Figure6: Benchmarks results (execution time is given in seconds).